

A statistical model for the estimation of natural gas consumption

Jiří Vondráček, Emil Pelikán, Ondřej Konár *, Jana Čermáková,
Kryštof Eben, Marek Malý, Marek Brabec

*Academy of Sciences of the Czech Republic, Institute of Computer Science, Pod Vodárenskou věží 2,
182 07, Praha 8, Czech Republic*

Available online 31 December 2007

Abstract

In this paper we present a statistical approach to natural gas consumption estimation of individual residential and small commercial customers. The approach is based on nonlinear regression principles. Parameters are estimated using mainly two real data sets – ordinary (approximately annual) meter readings of almost all customers and additional (approximately monthly) meter readings designed and operated within the frame of cooperation between the Institute of Computer Science of the Czech Academy of Sciences (ICS) and the West Bohemian Gas Distribution Company, a part of the RWE Group (WBG). The model was tested on various data sets. It has broad applicability in many areas of gas industry.

© 2007 Published by Elsevier Ltd.

Keywords: Nonlinear regression; Gas consumption modeling

1. Introduction

Natural gas distribution and trade companies need detailed information about their customers' behavior: about the consumption, response to weather conditions, seasonal character of natural gas usage etc. However, most gas distribution companies do not measure the consumptions of all their customers continuously due to technical and/or economical reasons. In the Czech Republic, meter readings are taken once a year in the case of residential and small commercial customers. Furthermore, each customer is measured in different time depending on the meter-reading itinerary. (Consequently, only a part of all customers – roughly one-twelfth – is billed during a month.) Therefore, gas distribution companies are looking for mathematical models which will be able to estimate gas consumption of individual customers for a given time period (e.g. days, months) with sufficient accuracy. Modeling of natural gas consumption of the whole segments of customers (e.g. cooking, water heating, space heating segments, etc.) are also required. Such models have broad applicability in areas such as unbilled revenues estimation, price elasticity estimation, temperature gradient

* Corresponding author. Tel.: +420 266 053 240.

E-mail address: konar@cs.cas.cz (O. Konár).

estimation, analysis of migrations between price segments, capacity pool control and optimization, and natural gas consumption forecasting.

Statistical models for natural gas consumption have been developed since the 1960's [1]. Various approaches have been published – a statistical model based on conditional demand analysis of end-used gas demand [2], an artificial neural network approach [3], a natural gas market equilibrium models for the US market [4] and other models [5–7]. However, many other publications or studies concerning mathematical modeling of natural gas consumption are not readily available (compared with the situation in electric power industry). The reason might be that “there is relatively little independent research on natural gas being carried out anywhere in the world” (as written in [8]).

In this paper we present a statistical model which was primarily developed for estimation of unbilled supply of natural gas, but it is applicable for solving many other tasks [9]. The model was developed in the frame of cooperation between the Institute of Computer Science (ICS) and the West Bohemian Gas distribution Company (WBG) in the Czech Republic [10,11]. The research part of the work was financed by the Grant Agency of the Czech Academy of Sciences.

2. The statistical model description

The ICS research group in cooperation with WBG developed a statistical model for natural gas consumption estimation. It is a nonlinear regression model with individual customer effect, typical time-dynamics part and the temperature correction.

2.1. Customer segmentation

Consumption characteristics, mainly the shape of seasonality pattern and strength of temperature dependence depend on the natural gas end-use. On that account the customers are classified into several segments for modeling purposes. Parameters of the model are then estimated separately for each customer segment. Two criteria are used for customer classification:

- (1) *Client type* (residential or commercial customers) which is more or less related to a consumption level.
- (2) *Natural gas end-use* – we consider combinations of the following end-uses: cooking, water heating, space heating and technology (which is relevant only for commercial customers).

This way we get seven segments for residential customers (all combinations of first three end-uses) and nine segments for commercial customers (two additional technological segments – including and excluding space heating). The total number of customer segments is 16 then (see Table 1).

2.2. Estimation of a particular customer's consumption

For consumption $C_{ik}(t)$ of a customer i from segment k in a day t we consider model

$$C_{ik}(t) = \mu_{ik}\Phi_k(t) + \varepsilon_{ikt}, \quad (1)$$

where μ_{ik} is an individual effect of a particular customer i from segment k , function $\Phi_k(t)$ of a day t is a typical time-dynamics of customers from segment k and ε_{ikt} is a random error, which is supposed to have a zero mean and variance proportional to the term $\mu_{ik}\Phi_k(t)$. The last assumption reflects the observed fact that consumption variability grows with the consumption level.

The model according to Eq. (1) is additive, consequently for a time period τ (e.g. an interval $[t_1, t_2]$) it holds

$$C_{ik}(\tau) = \mu_{ik} \sum_{t \in \tau} \Phi_k(t) + \sum_{t \in \tau} \varepsilon_{ikt}. \quad (2)$$

The typical time-dynamics function $\Phi_k(t)$ has the following form:

$$\Phi_k(t) = \Psi(t)e^{-\gamma_k f(T_t, N_t)} + p_k, \quad (3)$$

Table 1
The relative error of the estimate in various customer segments

Client type	Gas end-use	Monthly read		Annual read	
		Number of customers	Error (%)	Number of customers	Error (%)
Residential	S	143	1.08	9737	0.55
Residential	W	–	–	245	0.01
Residential	S + W	151	1.67	10495	0.00
Residential	C	45	0.48	98761	0.03
Residential	S + C	247	0.25	24688	0.36
Residential	W + C	16	5.60	9988	0.03
Residential	S + W + C	277	1.24	39838	0.24
Commercial	S	226	1.20	5836	0.11
Commercial	W	–	–	177	0.04
Commercial	S + W	115	0.81	2382	1.87
Commercial	C	13	8.34	492	5.02
Commercial	S + C	50	2.46	831	1.50
Commercial	W + C	–	–	62	0.12
Commercial	S + W + C	42	0.67	691	3.11
Commercial	T	20	5.36	501	4.14
Commercial	T + S	30	0.32	820	0.60
Total		1386	1.15	205544	0.30

S, space heating; W, water heating; C, cooking; T, technological use.

where $\Psi(t)$ is the seasonal part, γ_k is the segment-specific temperature dependence parameter, p_k is a constant (time independent) parameter for segment k , T_t is the actual (daily average) outdoor temperature in a day t (in °C) and N_t is a long-term normal (climatologic) temperature in a day t (in °C) and

$$f(T_t, N_t) = \tilde{T}_t - \tilde{N}_t, \quad (4)$$

where

$$\tilde{T}_t = \begin{cases} T_t & \text{if } T_t < 14 \\ 14 & \text{if } T_t \geq 14 \end{cases} \quad \tilde{N}_t = \begin{cases} N_t & \text{if } N_t < 14 \\ 14 & \text{if } N_t \geq 14 \end{cases} \quad (5)$$

Customers using natural gas for space heating exhibit a substantial difference in the consumption level between the summer and winter periods. For that reason, the modified form of function $\Phi_k(t)$ was introduced:

$$\Phi_k^{\text{mod}}(t) = \begin{cases} \Psi(t)e^{-\gamma_k f(T_t, N_t)} + p_k & \text{if } \bar{T}_t^{(3)} < 14 \\ q_k & \text{if } \bar{T}_t^{(3)} \geq 14 \end{cases} \quad (6)$$

where q_k is a segment-specific parameter interpreted as the mean summer consumption (summer consumption is supposed to be constant) and $\bar{T}_t^{(3)}$ is an average outdoor temperature of last three days (i.e. days $t, t-1, t-2$). The modified form $\Phi_k^{\text{mod}}(t)$ is then used for the “space heating segments” instead of $\Phi_k(t)$. For simplification just $\Phi_k(t)$ is written for all customer segments in following paragraphs.

At last, we denote the estimate of consumption $C_{ik}(t)$ based on the model according to Eq. (1) as

$$\hat{C}_{ik}(t) = \hat{\mu}_{ik} \hat{\Phi}_k(t). \quad (7)$$

2.3. Estimation of the individual effect

Individual effect μ_{ik} of a particular customer i in segment k is estimated by weighted least squares using data from ordinary (approximately annual) meter readings of WBG.

In fact the consumptions $C_{ik}(\tau_1), \dots, C_{ik}(\tau_{n_k})$ of a customer i in segment k in the time periods $\tau_1, \dots, \tau_{n_k}$ are known (from the last n_k meter readings). If the time periods are long enough, then we can assume that the

consumptions are stochastically uncorrelated. Parameter μ_{ik} is then estimated by the weighted least squares using the known consumptions of a particular customer:

$$\hat{\mu}_{ik} = \arg \min_{\mu} S^2(\mu) = \arg \min_{\mu} \sum_{j=1}^{n_k} \frac{(C_{ik}(\tau_j) - \mu \hat{\Phi}_k(\tau_j))^2}{\hat{\Phi}_k(\tau_j)} \quad (8)$$

where $\hat{\Phi}_k(\tau_j)$ is obtained by replacing parameters p_k , q_k and γ_k in Eq. (3) or in Eq. (6) with their estimates. In fact, we iterate between the individual effect and time-dynamics estimation steps in actual optimization procedure.

The solution of Eq. (8) leads to the estimate

$$\hat{\mu}_{ik} = \frac{\sum_{j=1}^{n_k} C_{ik}(\tau_j)}{\sum_{j=1}^{n_k} \hat{\Phi}_k(\tau_j)}. \quad (9)$$

Optimal number \hat{n}_k of meter readings used for estimation (the so-called *history depth parameter*) depends on segment k and is estimated in longer time intervals together with the other parameters of the model according to Eq. (1) (see Section 2.6).

If we substitute the estimate $\hat{\mu}_{ik}$ from Eq. (9) into Eq. (7) we get the consumption estimator:

$$\hat{C}_{ik}(t) = \hat{\mu}_{ik}(t) \hat{\Phi}_k(t) = \frac{\sum_{j=1}^{\hat{n}_k} C_{ik}(\tau_j)}{\sum_{j=1}^{\hat{n}_k} \hat{\Phi}_k(\tau_j)} \hat{\Phi}_k(t) = \frac{\hat{\Phi}_k(t)}{\sum_{j=1}^{\hat{n}_k} \hat{\Phi}_k(\tau_j)} \sum_{j=1}^{\hat{n}_k} C_{ik}(\tau_j) \quad (10)$$

where $\hat{\Phi}_k(t) = \Psi(t)e^{-\hat{\gamma}_k f(T_t, N_t)} + \hat{p}_k$ (and accordingly for $\hat{\Phi}_k^{\text{mod}}(t)$). The estimate is a product of the data part $\sum_{j=1}^{\hat{n}_k} C_{ik}(\tau_j)$ and the model part $\frac{\hat{\Phi}_k(t)}{\sum_{j=1}^{\hat{n}_k} \hat{\Phi}_k(\tau_j)}$.

The model part behavior with $\hat{n}_k = 1$, 365 days long period and $\hat{\gamma}_k = 0$ is shown in Fig. 1. Two curves for space heating ($\hat{p}_k = 4$) and cooking ($\hat{p}_k = 2000$) segments are drawn. The curves represent a distribution of the annual consumption into particular days of a year. It is apparent that the cooking segment curve is much flatter than the space heating segment one. In fact, it is almost constant.

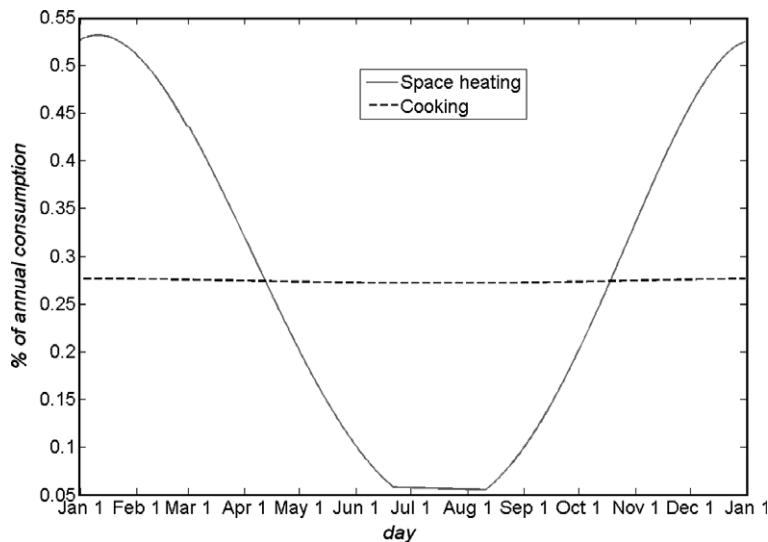


Fig. 1. Distribution of the annual consumption into particular days of a year for cooking and space heating segment.

2.4. The typical time-dynamics function

In this paragraph, we will derive the form of the typical time-dynamics function $\Phi_k(t)$ in Eq. (3). First we assume for simplicity that the actual outdoor temperature always equals to the long-term so-called “normal” temperature. Then there is no need for any temperature correction. Secondly, we will work just with the systematic part of the model, so we will omit the random term ε_{ikt} .

Let us denote $\psi(t)$ the consumption of all customers in a day t . The principal idea is that the consumption is splits into two parts – time independent part c and time dependent part $(\psi(t) - c)$. Consumption is distributed into these two parts differently in different segments. Therefore the total consumption $\phi_k(t)$ of segment k in a day t can be expressed as

$$\phi_k(t) = w_{k1}(\psi(t) - c) + w_{k2}c, \quad (11)$$

where w_{k1} and w_{k2} are the segment-specific weights. Customer i then consumes in a day t

$$C_{ik}(t) = \mu_{ik}\phi_k(t). \quad (12)$$

Similarly, as in Paragraph 2.3 we get an estimate of μ_{ik} :

$$\hat{\mu}_{ik} = \frac{\sum_{j=1}^{n_k} C_{ik}(\tau_j)}{\sum_{j=1}^{n_k} \hat{\phi}_k(\tau_j)}. \quad (13)$$

Then the estimate $\hat{C}_{ik}(t)$ of the consumption $C_{ik}(t)$ has a form of

$$\hat{C}_{ik}(t) = \hat{\mu}_{ik}\hat{\phi}_k(t) = \frac{\hat{\phi}_k(t)}{\sum_{j=1}^{n_k} \hat{\phi}_k(\tau_j)} \sum_{j=1}^{n_k} C_{ik}(\tau_j). \quad (14)$$

If we replace the term $\hat{\phi}_k(t)$ with $\eta\hat{\phi}_k(t)$ (and accordingly $\hat{\phi}_k(\tau_j)$ with $\eta\hat{\phi}_k(\tau_j)$) where η is an arbitrary constant, we will get the same estimate $\hat{C}_{ik}(t)$.

Hence if we multiply the term $\phi_k(t)$ in Eq. (12) by any constant, we get a model that gives the same estimate of consumption as the model according to Eq. (12). With a proper choice of a constant η we get the model

$$C_{ik}(t) = \mu_{ik}(\Psi(t) + p_k), \quad (15)$$

where $\Psi(t) = \frac{\psi(t)}{c}$ and p_k is a function of the weights w_{1k} , w_{2k} . The model according to Eq. (15) has simpler structure in contrast to the model according to Eq. (12) and it is more convenient for the implementation and parameter estimation process. Both models give the identical estimate of consumption. The term $\Psi(t)$ represents a multiplicative seasonality of a part of consumption above the minimal consumption p_k independent from the temperature effect.

Finally, we have to return to the real conditions where we need a temperature correction. We may expect that the time independent part of consumption p_k is also temperature independent. For a temperature correction we multiply the seasonal part $\Psi(t)$ by the temperature correction term $e^{-\gamma_k f(T_t, N_t)}$ and we get the formula Eq. (3).

2.5. Estimation of the seasonal part

It is a difficult task to estimate the seasonal part $\Psi(t)$ of the typical time-dynamics function. The reason is that there is a lack of information about seasonality in the annual meter readings. So there is a need of additional data. We used information about the total consumption of WBG residential and small commercial customers. Such consumption can be estimated from the total consumption of the whole company subtracting the measured consumption values of (large) commercial customers (unknown network losses are not taken into account). Since the consumption of some (large) commercial customers are measured monthly, the total daily consumption of the residential and small commercial customers had to be interpolated. The fifth-order polynomial was fitted on the monthly values using ordinary least squares method. Since values of monthly consumption of the particular segments are not available, common seasonal part $\Psi(t)$ is used for all customer segments.

The common seasonal part $\Psi(t)$ was estimated in the beginning of the project and remains fixed due to the lack of reasonable data for re-estimation. For that reason we write $\Psi(t)$ instead of more precise $\hat{\Psi}(t)$ and regard it as a fixed vector rather than a parameter.

2.6. Estimation of the remaining parameters from the annual meter readings

The “history depth parameter” n_k , temperature dependence rate γ_k , time independent part of consumption p_k and summer (average) consumption q_k (for the “space heating” segments) are estimated using annual meter readings of all WBG small commercial and residential customers. Customers with unstable and nonstandard behavior (found by the WBG experts) and those with more than one gas-meter are excluded.

For each segment k the parameter estimates are computed by minimizing the function:

$$K(p_k, q_k, \gamma_k, n_k) = \frac{1}{N_k^*} \sum_{i \in k} \frac{|C_{ik}(\tau_{ik}) - \hat{C}_{ik}(\tau_{ik})|}{\|\tau_{ik}\|} \quad (16)$$

where $C_{ik}(\tau_{ik})$ is the last measured consumption of the customer i from the segment k in a time period τ_{ik} , $\hat{C}_{ik}(\tau_{ik})$ is the estimate according to Eq. (7) of consumption in that period with parameters μ_{ik} , γ_k , p_k , q_k and n_k replaced by their estimates; parameter μ_{ik} is estimated as described in 2.3, $\Psi(t)$ is replaced by the estimate described in 2.5. $\|\tau_{ik}\|$ is the number of days of the period τ_{ik} and N_k^* is the total number of customers in segment k .

Optimization is currently performed on a grid of parameter values by the gradient descent method.

2.7. Role of additional data sets

As we have described in the previous paragraph, regular meter readings of (almost) all WBG small customers are used for parameter estimation. The average time period length $\|\tau_{ik}\|$ is approximately one year according to the itinerary, so the amount of information about monthly or daily consumption is rather low. This affects the quality of parameter estimates.

In order to improve the performance of the developed model, ICS in cooperation with WBG has been collecting additional monthly data of about 1700 customers since September 2004. These data have been used for parameter estimation recently. We tested these two approaches (using ordinary approximately annual meter readings and additional approximately monthly readings) on testing data set formed by daily total consumptions of 1055 customers in a selected Western Bohemian city. In Fig. 2 we can see some improvement especially during the summer period (low consumption) after including the additional data into the parameter

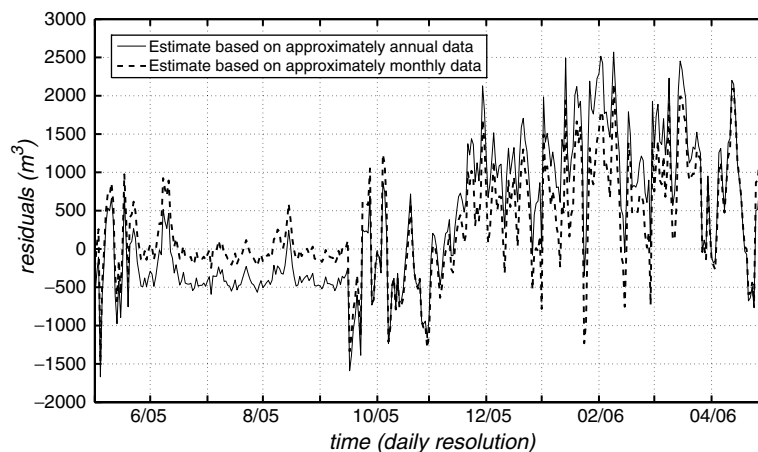


Fig. 2. Residuals from the estimates based on two different datasets – approximately annual and approximately monthly meter readings, period May 2005–May 2006.

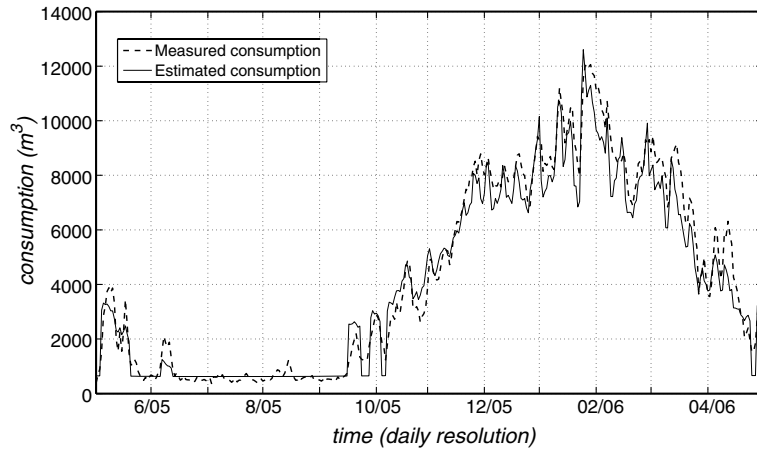


Fig. 3. Comparison of estimated (using approximately monthly meter readings) and measured total consumption of an experiment locality in Western Bohemia, period May 2005–May 2006.

estimation process. Comparison of estimated (using approximately monthly meter readings) and measured total consumption of the corresponding testing locality is shown in Fig. 3.

The parameter estimation differs a little bit due to the character of additional monthly data. Parameters p_k , q_k and γ_k were estimated by minimizing the function

$$\tilde{K}(p_k, q_k, \gamma_k) = \frac{1}{N_k^*} \sum_{i \in k} \sum_{j=1}^{v_{ik}} \frac{|C_{ik}(\tau_{ikj}) - \hat{C}_{ik}(\tau_{ikj})|}{\|\tau_{ikj}\|} \quad (17)$$

where $C_{ik}(\tau_{ik1}), \dots, C_{ik}(\tau_{ikv_{ik}})$ are the measured monthly consumptions of customer i from segment k , $\hat{C}_{ik}(\tau_{ik1}), \dots, \hat{C}_{ik}(\tau_{ikv_{ik}})$ are the estimates according to Eq. (7) of consumption in corresponding time periods with parameters μ_{ik} , γ_k , p_k , q_k and n_k replaced by their estimates; parameter μ_{ik} is estimated as described in 2.3, $\Psi(t)$ is replaced by estimate described in 2.5. $\|\tau_{ikj}\|$ is the number of days of the period τ_{ikj} and N_k^* is the total number of customers in segment k . Optimization was performed on a grid of possible parameter values testing all the values in the grid. The optimal history depth parameter n_k was found using the ordinary (approximately annual) meter readings and the estimates $\hat{p}_k, \hat{q}_k, \hat{\gamma}_k$ minimizing function $K(p_k, q_k, \gamma_k)$ from Eq. (17).

3. The model performance

The relative error of the estimate

$$\frac{|\hat{C}_{ik}(\tau) - C_{ik}(\tau)|}{C_{ik}(\tau)} \quad (18)$$

is acceptably low if we estimate the consumption in longer time periods (e.g. 1 year). In Table 1 we show the relative estimation errors for all customer segments. The errors were computed on two data sets. The first set was the sample of monthly read customers. The estimated consumption in a 1.5 year time interval was compared to their measured consumption (sum of the monthly consumptions). An outlier analysis was performed on the additional monthly meter readings and the outliers were excluded. It is apparent that the error is inversely proportional to the segment frequency (number of customers in a segment). Segments with less than 10 customers were omitted.

The second sample was much larger. The last regular meter reading of each customer was compared with the estimated consumption in a corresponding time period. The period was different for each customer and was in average about 1 year long. It is apparent that the errors are lower in the case of larger sample. The

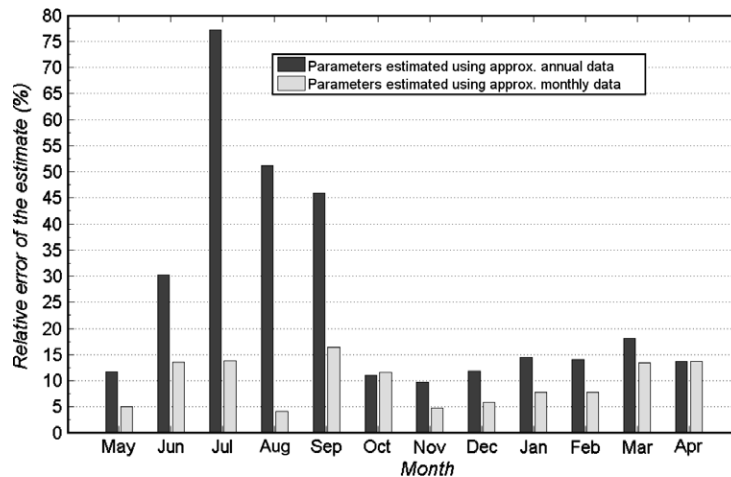


Fig. 4. Relative error of two estimates (using approximately annual and approximately monthly data) of total consumption in experimental locality in West Bohemia for month periods from May 2005 to April 2006.

error is large in commercial segments due to smaller number of customers and due to more variable consumption behavior.

The relative error of the estimate according to Eq. (18) for shorter time periods (e.g. a month) is larger. With the parameter estimates based on the ordinary meter readings (taken in more or less yearly intervals), the error according to Eq. (18) grows over 50% in summer months. When we use the additional monthly data, the error according to Eq. (18) decreases significantly as it is evident from Fig. 4.

4. Future challenges

The proposed model forms a core of the operational software system which used by WBG. The parameters are re-estimated every 3 months.

The right choice of optimization methods is very important. The first challenge is the proper choice of a criterion function. That is partially influenced by the needs of WBG. In a further development we will try to use one common criterion function for estimating all the parameters (which is not true now).

Next challenge is the optimization algorithm. Since the form of criterion function (the currently applied form or, e.g., weighted least squares) is complicated, there is no simple formula for the estimator like the ordinary least squares (OLS) estimator in linear regression. Thus various numeric methods are tested in order to improve the running and the results of the optimization process.

The seasonal part $\Psi(t)$ of the model base $\Phi_k(t)$ does not correspond to the real behavior of the consumption in some segments. Hence we examine the possibilities of estimating this function from the additional (approximately monthly) data with higher precision using a proper parametrization or nonparametric approach. Furthermore, we test various nonlinear forms of function f from the temperature correction term in Eq. (4).

5. Conclusions

A nonlinear regression model for estimation of natural gas consumption was presented. The development of the model is a long-term process which will run with a grant support at least until 2009. In this paper we presented the current status of development. Some challenges for further research were also indicated.

The main objective of this paper was rather than introducing a new statistical method or discovery making a view into the work of application oriented researcher who has to apply and modify known common methods and is limited by the quality and quantity of data on one hand and the demands of the end user on the other hand.

Acknowledgements

This paper was supported by the GA AS CR Grant No. 1ET400300513 and by the institutional research plan of ICS AS CR No. AV0Z10300504.

The authors thank the West Bohemian Gas Distribution Company (namely Jaroslav Matějovic, Olga Volfová–Naxerová, Josef Bečvář and Ladislava Blahová) for a valuable cooperation.

References

- [1] Balestra P, Nerlove M. Pooling cross section and time series data in the estimation of a dynamic model: the demand for natural gas. *Econometrica* 1966;34(3):585–612.
- [2] Bartels R, Fiebig DG, Nahm D. Regional end-use gas demand in Australia. *Econ Rec* 1996;72(219):319–31.
- [3] Suykens J, Lemmerling Ph, Favoreel W, De Moor B, Crepel M, Briol P. Modelling the Belgian gas consumption using neural networks. *Neural Process Lett* 1996;4(3):157–66.
- [4] Gabriel SA, Kiet S, Zhuang J. A mixed complementarity-based equilibrium model of natural gas markets. *Oper Res* 2005;53(5):799–818.
- [5] Gümrah F, Katircioglu D, Aykan Y, Okumus S, Kiliñçer N. Modeling of gas demand using degree-day concept: case study for Ankara. *Energy Source* 2001;23:101–14.
- [6] Gutiérrez R, Nafidi A, Gutiérrez Sánchez R. Forecasting total natural-gas consumption in Spain by using the stochastic Gompertz innovation diffusion model. *Appl Energy* 2005;80:115–24.
- [7] Mackay RM, Probert SD. Forecasting the United Kingdom's supplies and demands for fluid fossil-fuels. *Appl Energy* 2001;69:161–89.
- [8] Natural Gas Research Program web pages, Oxford Institute for Energy. URL: www.oxfordenergy.org.
- [9] Čermáková J, Bečvář J, Naxerová O, Brabec M, Brabec T, Konár O, et al. Natural gas consumption modeling: customers without course measurement. In: *Proceedings of the seventh SIMONE Congress (CD)*, Lednice; October 11–14, 2005.
- [10] Čermáková J, Matějovic J, Naxerová O, Bečvář J, Brabec M, Brabec T, et al. Mathematical modeling of natural gas consumption without continuous measurement. *Plyn* 2005;2:34–7 (in Czech).
- [11] Čermáková J, Volfová–Naxerová O, Bečvář J. The determination of the volume of the uninvoiced gas using the virtual invoicing and the GAMMA mathematical model. *Plyn* 2005;4:76–9 (in Czech).